

Collezioni digitali di periodici in Italia e in Europa: standard, applicazioni, valutazioni, prospettive

Arco (TN) - 15-16 Novembre 2007

Come contare gli utenti? Le basi tecnologiche per rilevare gli utenti di un sito

Zeno Tajoli – tajoli@cilea.it



Abstract

In questo workshop si vogliono illustrare le basi e le problematiche connesse alla rilevazione degli utenti di un sito, le loro caratteristiche, le operazioni che compiono. Per compiere queste rilevazioni si illustrerà come funziona in generale la comunicazione tra un utente e un sito e come gli utenti siano anonimi, ma lascino sempre dei dati, che una volta interpretati, possono dire molto di loro. In particolare verranno illustrate le problematiche sottese a queste operazioni, per permettere al bibliotecario di capire meglio cosa significano esattamente i dati forniti dai software che valutano l'uso di risorse digitali accessibili a tutti gli utenti di Internet.

Prerequisiti

- ☐ Uso di Internet
- ☐ Generica conoscenza dei PC
- ☐ Interesse a guardare dietro le quinte
- ☐ Nessun vero background tecnico richiesto
- ☐ Si parte dalla base

Punti salienti

- ☐ Come si comunica?
- ☐ Che informazioni ho?
- ☐ Come estrarre i dati.
- ☐ Superare i limiti presentati
- ☐ Un esempio

Come si comunica ?

Sono necessari meccanismi software per permettere ai vari computer di dialogare di gestire la comunicazione

- ☐ protocolli (convenzioni) di comunicazione
- ☐ meccanismi di indirizzamento (come identificare un computer)
- ☐ spedizione sulle connessioni opportune

Come si comunica ?

- ☐ invio e ricezione di messaggi
- ☐ verifica correttezza dei messaggi durante la trasmissione
- ☐ protezione dei messaggi (per evitare intercettazione)
- ☐ ottimizzazione della comunicazione
- ☐ gestione del traffico sulla rete

Come si comunica ?

Un protocollo umano e un protocollo di reti di computer:

- 1) Ciao
- 2) Ciao
- 3) Hai l'ora?
- 4) 2:00

- a) TCP connection
- b) request
- c) TCP connection
- d) reply.
- e) Get <http://www.di.unito.it/index.htm>
- f) <file>

Domanda: Altri protocolli umani?

Come si comunica ?

Protocolli umani:

- "Che ora è?"
- "Ho una domanda"
- Presentazioni...

... messaggi specifici vengono spediti

... azioni specifiche sono compiute quando i messaggi sono ricevuti, o in seguito ad altri eventi

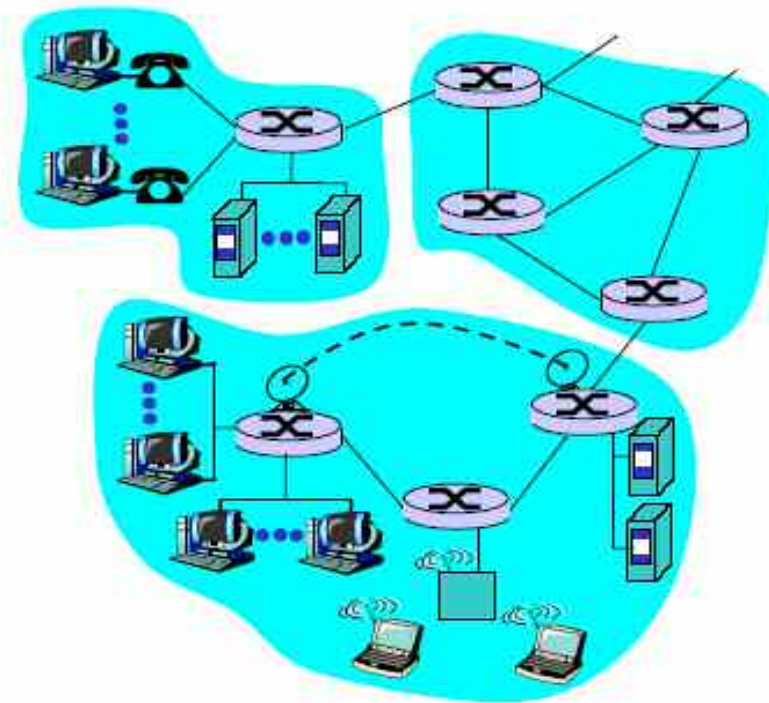
Protocolli di rete:

- macchine invece di esseri umani
- Tutte le attività di comunicazione in Internet sono governate da protocolli

I protocolli definiscono formato e ordine dei messaggi spediti e ricevuti tra entità della rete, e le azioni da compiere in seguito alla ricezione e/o trasmissione dei messaggi o di altri eventi

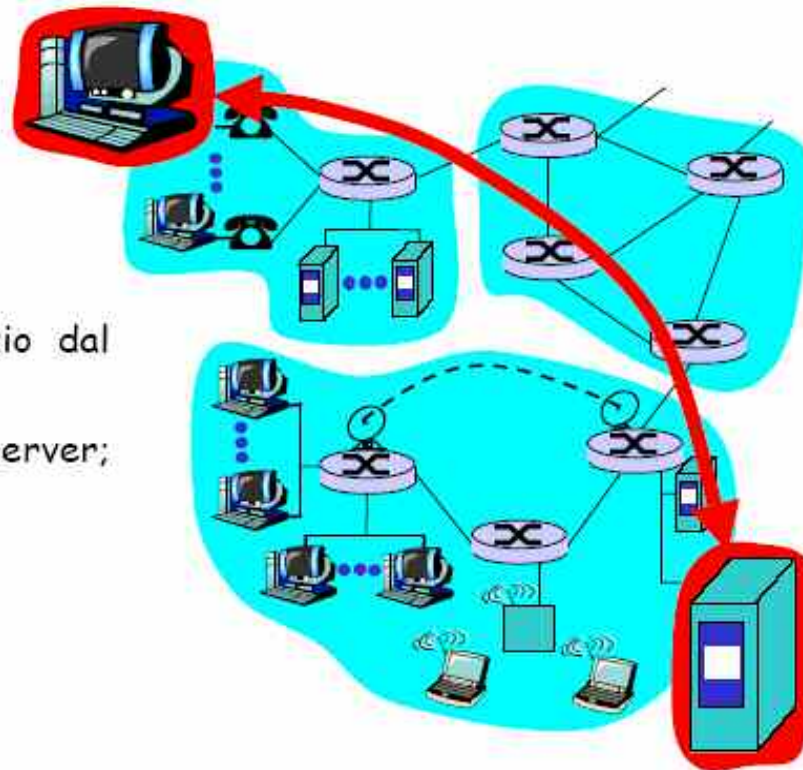
Come si comunica ?

- **network edge:** applicazioni ed host
- **network core:**
 - router
 - rete di reti
- **reti di accesso, mezzi trasmissivi:** canali di comunicazione



Come si comunica ?

- **end systems (host):**
 - Eseguono programmi applicativi
 - e.g., WWW, email
 - al "bordo della rete"
- **modello client/server**
 - il client richiede, riceve servizio dal server
 - e.g., WWW client (browser)/ server; email client/server
- **modello peer-peer:**
 - interazione tra host simmetrica
 - e.g.: Gnutella, KaZaA



Come si comunica ?

Obiettivo: trasferimento dati tra host

- *handshaking*: fase di preparazione antecedente al trasferimento dati
 - Ciao - Ciao nel protocollo umano
 - *Stabilire uno "stato"* nei due host comunicanti
- TCP - Transmission Control Protocol
 - Servizio di scambio dati di tipo connection-oriented di Internet

Servizio TCP [RFC 793]

- *Trasferimento affidabile ed ordinato di byte di un flusso dati*
 - perdite: conferma di ricezione (acknowledgement) e ri-trasmissione
- *Controllo di flusso*
 - Il mittente non sovraccaricherà il ricevitore
- *Controllo di congestione:*
 - I mittenti diminuiscono la loro velocità di spedizione quando la rete si congestiona

Come si comunica ?

Obiettivo: trasferimento dati tra host

- Esattamente lo stesso!
- **UDP** - User Datagram Protocol [RFC 768]: Servizio connectionless di Internet
 - Senza handshaking
 - Trasferimento dati non-affidabile
 - senza controllo di flusso
 - senza controllo congestione

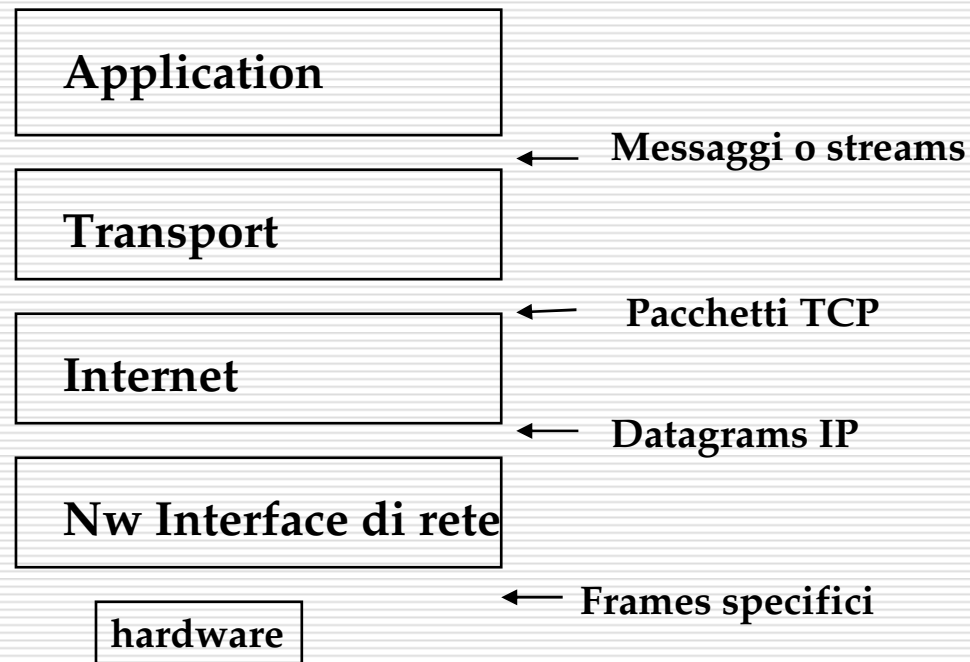
Applicazioni che usano TCP:

- HTTP (WWW), FTP (trasferimento file), Telnet (login remoto), SMTP (email)

Applicazioni che usano UDP:

- streaming media, teleconferencing, Internet telephony

Come si comunica ?



Il www in realtà usa due protocolli, il TCP e l'IP

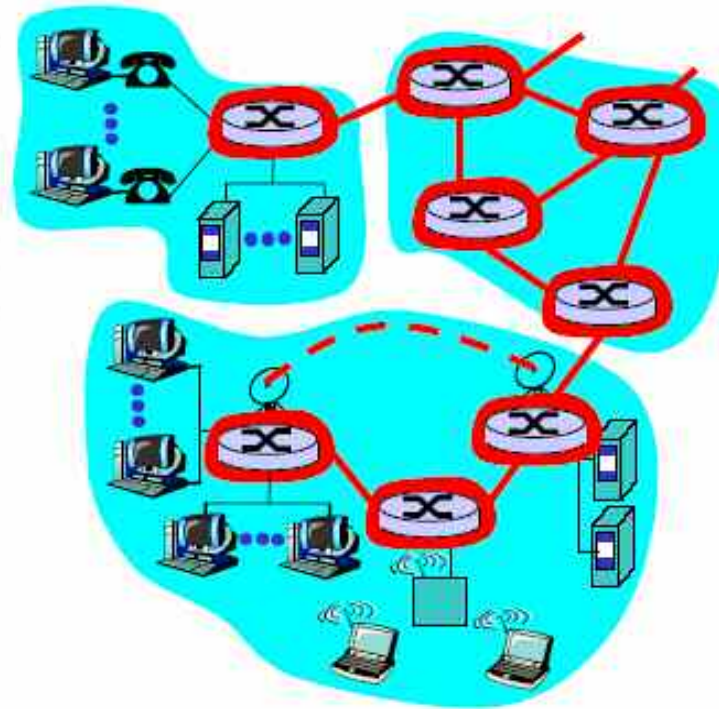
Come si comunica ?

Il protocollo IP

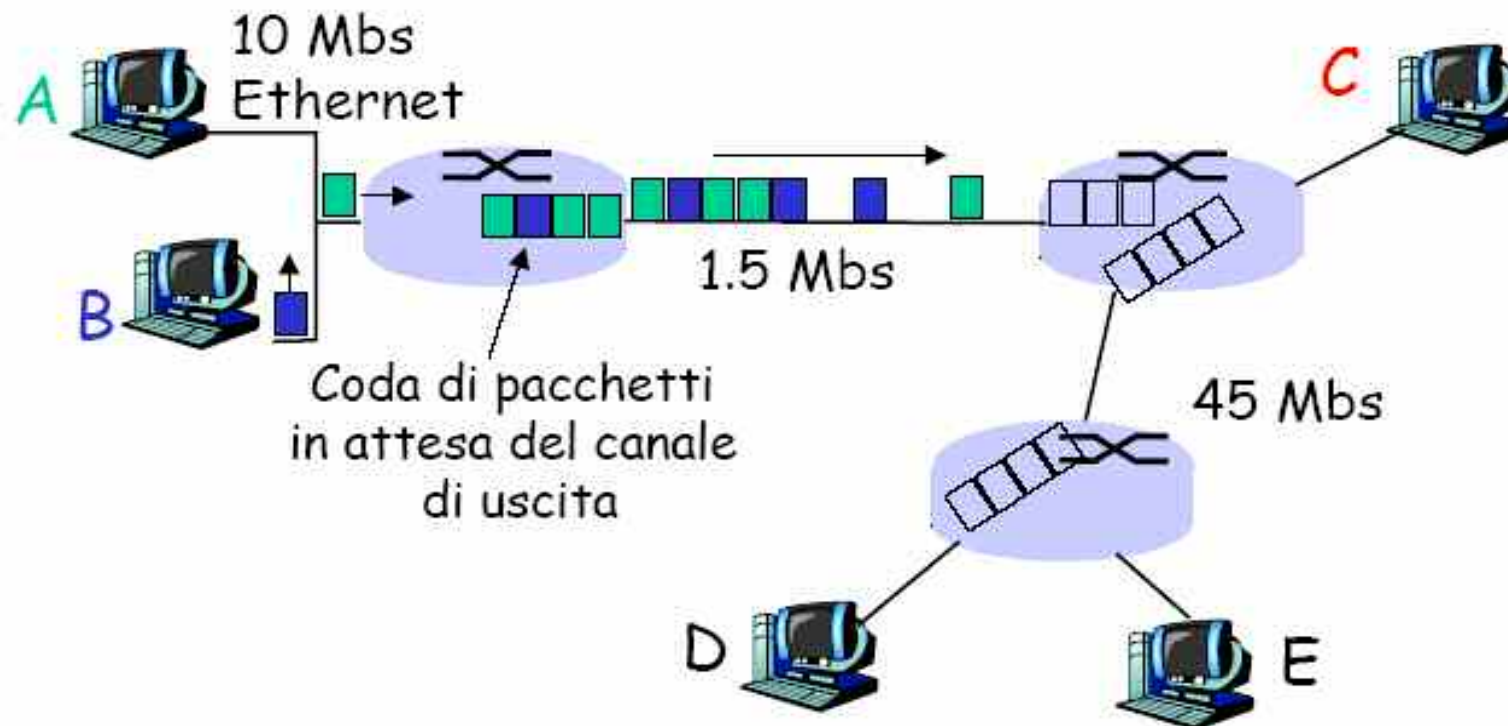
- protocollo che lavora con pacchetti (datagram) indipendenti e di formato predefinito
- i nodi possono eliminare datagram se mancano le risorse
- nessun algoritmo di correzione sui datagram
- frammentazione lungo la strada e riassemblaggio gestito alla destinazione
- nessuna garanzia di instradamento dei datagram
- riconfigurazione automatica della strada percorsa in caso di malfunzionamenti
- distruzione dei datagram se ripetitivi

Come si comunica ?

- Maglia di router interconnessi
- Domanda fondamentale: come vengono trasferiti i dati attraverso la rete?
 - Commutazione di pacchetto: i dati sono spediti attraverso la rete in quantità discrete chiamate **pacchetti**



Come si comunica ?



Come si comunica ?

Ogni flusso dati viene diviso in
pacchetti

- I pacchetti degli utenti A e B *condividono* risorse di rete
- Ogni pacchetto usa tutta la larghezza di banda (capacità di trasmissione in bit al secondo) del canale
- Risorse usate quando sono necessarie

Contesa delle risorse:

- La richiesta aggregata di risorse può eccedere l'ammontare disponibile
- congestione: i pacchetti si accodano ed attendono l'uso del canale
- store and forward: pacchetti ricevuti interamente prima di essere spediti

Come si comunica ?

- Obiettivo: spostare pacchetti tra router, dal host sorgente all' host destinatario
- **Caratteristiche**:
 - *L'indirizzo destinazione* determina il prossimo passo
 - Le strade (route) possono variare durante le sessioni
 - I router NON mantengono informazioni sullo stato delle connessioni

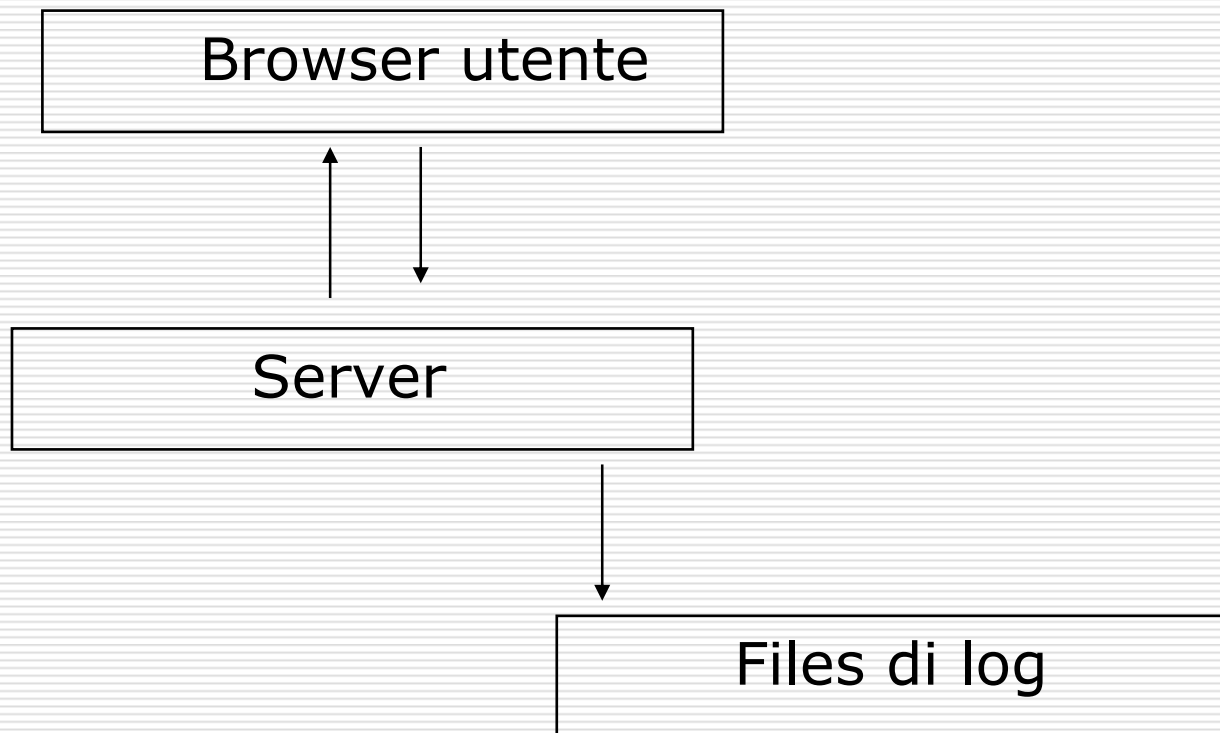
Come si comunica ?

- ☐ Una struttura di comunicazione complessa
- ☐ Dimensionabile da 2 a n sistemi
- ☐ Senza preconoscenza di tutta la struttura
- ☐ Senza totale anonimato
- ☐ Senza totale sicurezza di chi è e cosa ha realmente ricevuto

Come si comunica ?

- ❑ La comunicazione è tra macchine
- ❑ Il dato fondamentale è l' IP
- ❑ Ma non tutti l'hanno fisso
- ❑ I grandi fornitori di connettività (TIM, AOL, Tiscali, etc.) ne hanno un pool che fanno usare, girandoli, tra gli utenti finali
- ❑ Anche server con IP dinamico ([Dynamic DNS](#))

Che informazioni ho ?



Che informazioni ho ?

- ❑ Via via che la comunicazione client/server va avanti, il server scrive quanto avviene su un file
- ❑ Di norma i file sono due,
 - Il primo per le operazioni normali
 - Il secondo per quanto ha causato errori
- ❑ Si analizzerà quanto fa il server web Apache 2, simile é il comportamento di IIS di Microsoft

Che informazioni ho ?

- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /archive/00011300/ HTTP/1.1" 200 12384 "http://www.google.pt/search?q=caf&hl=pt-PT&lr=lang_pt&start=170&sa=N" "Mozilla/4.0(MSIE 6.0)"
- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /eprints.css HTTP/1.1" 200 3723 "http://eprints.rclis.org/archive/00011300/" " Mozilla/4.0(MSIE 6.0)" "
- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /images/imatge.jpg HTTP/1.1" 200 38180 "http://eprints.rclis.org/archive/00011300/" " Mozilla/4.0(MSIE 6.0)"
- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /images/mon.jpg HTTP/1.1" 200 14692 "http://eprints.rclis.org/archive/00011300/" "Mozilla/4.0(MSIE 6.0)"
- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /images/dibuix.jpg HTTP/1.1" 200 781 "http://eprints.rclis.org/archive/00011300/" "Mozilla/4.0(MSIE 6.0)"
- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /images/reflink.png HTTP/1.1" 200 353 "http://eprints.rclis.org/archive/00011300/" "Mozilla/4.0(MSIE 6.0)"
- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /email.js HTTP/1.1" 200 194 "http://eprints.rclis.org/archive/00011300/" " Mozilla/4.0(MSIE 6.0)"
- ❑ 72.226.75.yy "-" "-" [11/Nov/2007:06:43:14 +0100] "GET /7876/ HTTP/1.1" 200 5 "-" "Mozilla/5.0 Firefox/2.0.0.9"

Che informazioni ho ?

- ☐ IP-address remoto
- ☐ IP-address locale
- ☐ Dimensione della risposta (senza gli header)
- ☐ Dati del cookie inviato dal browser
- ☐ Tempo per gestire la richiesta sul server
- ☐ Variabili interne del server
- ☐ Nome del file inviato
- ☐ Protocollo di richiesta
- ☐ Informazioni che il browser dà di se stesso
 - Che lingua richiede
 - Che set di caratteri usa
 - Che pagina ha visto prima
 - Come si chiama il software del browser
 - Altre informazioni molto variabili tra i diversi software
- ☐ LogNome dell'utente remoto se autenticato
- ☐ Il modo di fare la richiesta

Che informazioni ho ?

- ☐ Le intestazioni inviate nella risposta
- ☐ La porta di comunicazione
- ☐ I numeri interni del servizio nella memoria (PID)
- ☐ La parte domanda di una richiesta ricevuta
- ☐ La richiesta ricevuta
- ☐ Lo status della comunicazione
- ☐ L'ora in cui si è ricevuta la richiesta
- ☐ Nome associato al LogName dell'utente remoto se autenticato
- ☐ La richiesta ricevuta senza la eventuale parte domanda
- ☐ Nome del server
- ☐ Status della connessione dopo l'invio di quanto chiesto
- ☐ Bytes ricevuti tutto compreso
- ☐ Bytes inviati tutto compreso

Riferimento completi: RFC 2616 e modulo Apache 2 mod_log_config

Che informazioni ho ?

- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /archive/00011300/ HTTP/1.1" 200 12384 "http://www.google.pt/search?q=caf&hl=pt-PT&lr=lang_pt&start=170&sa=N" "Mozilla/4.0(MSIE 6.0)"

Da IP 213.58.139.xx un browser MSIE chiama eprints.rclis.org e chiede la home di 11300, viene da una query su google, la richiesta è andata buon fine e sono stati inviati 12384 bytes

- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /eprints.css HTTP/1.1" 200 3723 "http://eprints.rclis.org/archive/00011300/" "Mozilla/4.0(MSIE 6.0)"

Sempre allo stesso inviato con successo il file eprints.css

- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /images/imatge.jpg HTTP/1.1" 200 38180 "http://eprints.rclis.org/archive/00011300/" "Mozilla/4.0(MSIE 6.0)"

Sempre allo stesso inviato con successo il file imatge.jpg

- ❑ 213.58.139.xx "-" "-" [11/Nov/2007:06:43:13 +0100] "GET /images/mon.jpg HTTP/1.1" 200 14692 "http://eprints.rclis.org/archive/00011300/" "Mozilla/4.0(MSIE 6.0)"

Sempre allo stesso inviato con successo il file mon.jpg

- ❑ Vari invii di diversi files allo stesso IP

Che informazioni ho ?

- ❑ 72.226.75.yy "-" "-"
[11/Nov/2007:06:43:14 +0100] "GET
/7876/ HTTP/1.1" 200 5 "-" "Mozilla/5.0
Firefox/2.0.0.9"

Connessione con successo di un altro IP

- ❑ I log di esempio hanno questa sintassi:
- ❑ "%h \"%I\" \"%u\" %t \"%r\" %>s %b
\"%{Referer}i\" \"%{User-Agent}i\""
- ❑ I dati sono in stretto ordine di tempo
- ❑ Dunque si accavallano l'un con l'altro

Come estraggo i dati.

- ☐ Ma quando finisce la sessione ?
- ☐ Posso collegare i vari passi dell'utente ?
- ☐ Se ritorna lo stesso IP è la stessa macchina ?
- ☐ Stessa macchina = stessa persona ?

Come estraggo i dati.

- ❑ I raggruppamenti e le successive analisi si fanno sfruttando di base:
 - Il valore dell'IP
 - L'orario dell'operazione
- ❑ La fine sessione viene definita attraverso un intervallo temporale minimo tra una nuova evenienza dello stesso IP
- ❑ Un limite dice che è lo stesso utente che fa una nuova sessione
- ❑ Un secondo limite dice che è un utente diverso.
- ❑ Un esempio:
 - Tra 0 e 5 minuti: stesso utente e stessa sessione
 - Tra 5 minuti e 1 ora: stesso utente ma diversa sessione
 - Oltre 1 ora: diverso utente

Come estraggo i dati.

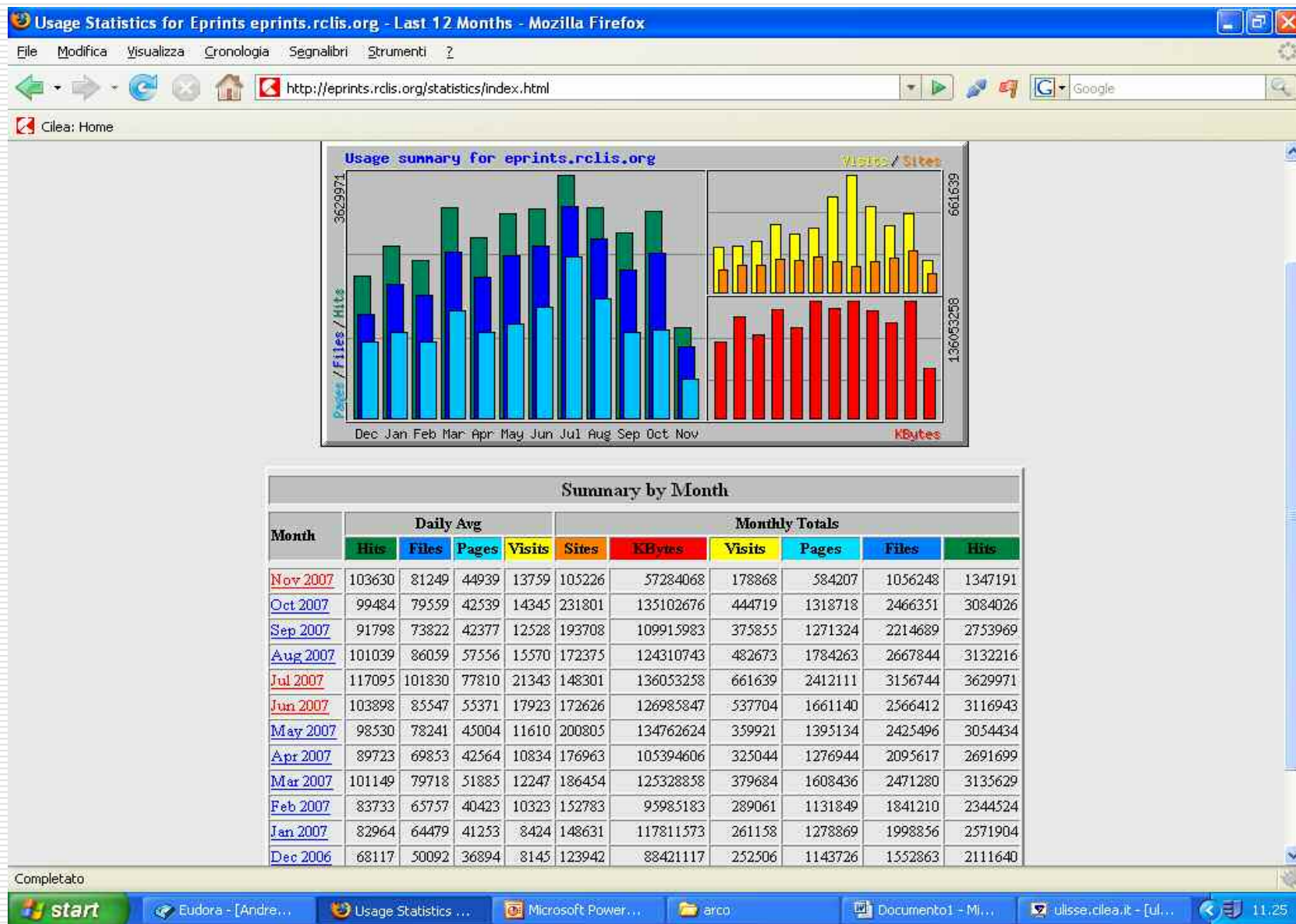
- ❑ La gestione dei limiti temporali è il passaggio cruciale le elaborazione.
- ❑ I dati sui files e su quanto è stato inviato sono dunque precisi e certi
- ❑ I dati sugli utenti sono dunque delle stime.
- ❑ Diversi software usano diversi algoritmi per queste stime

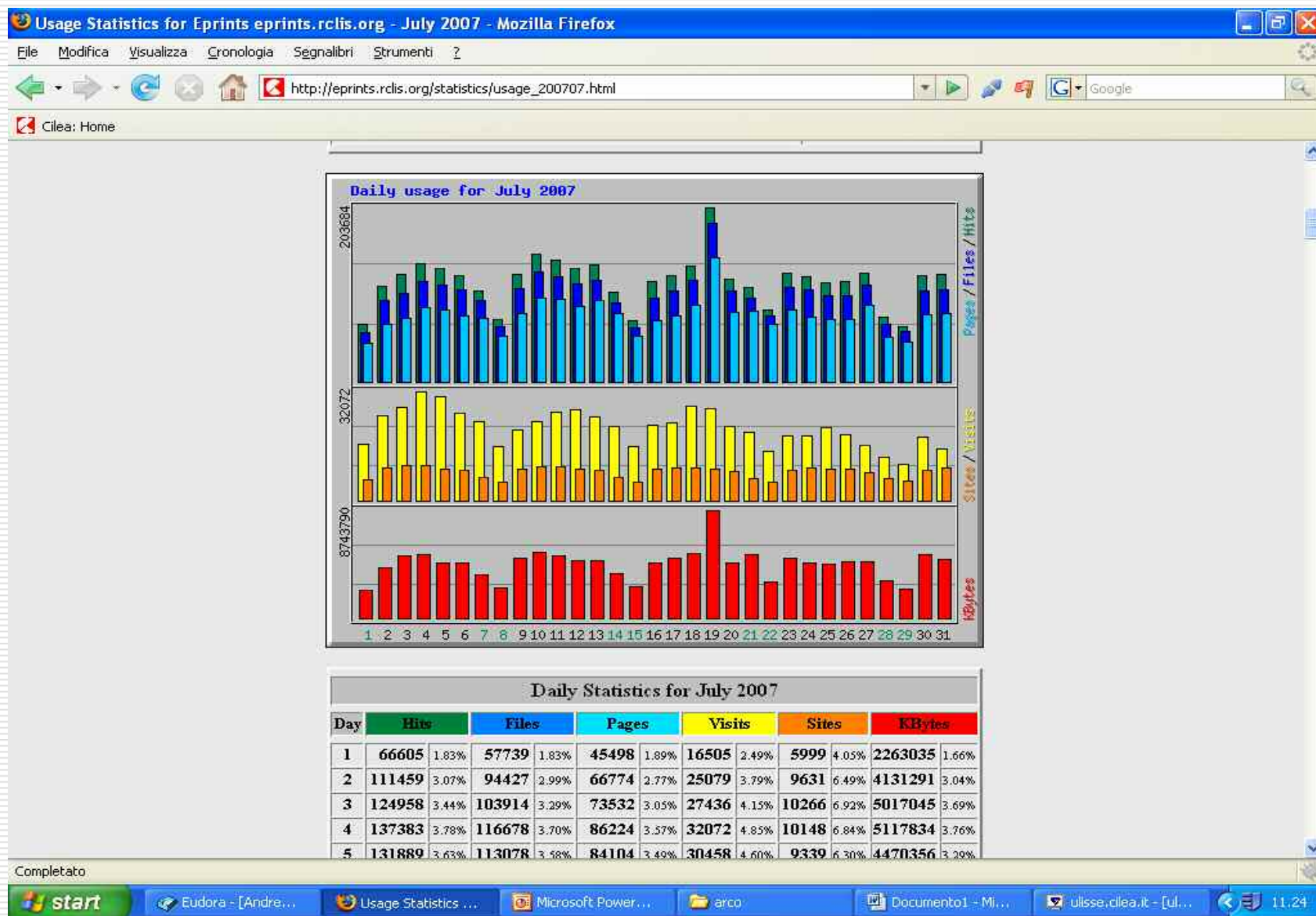
Come estraggo i dati.

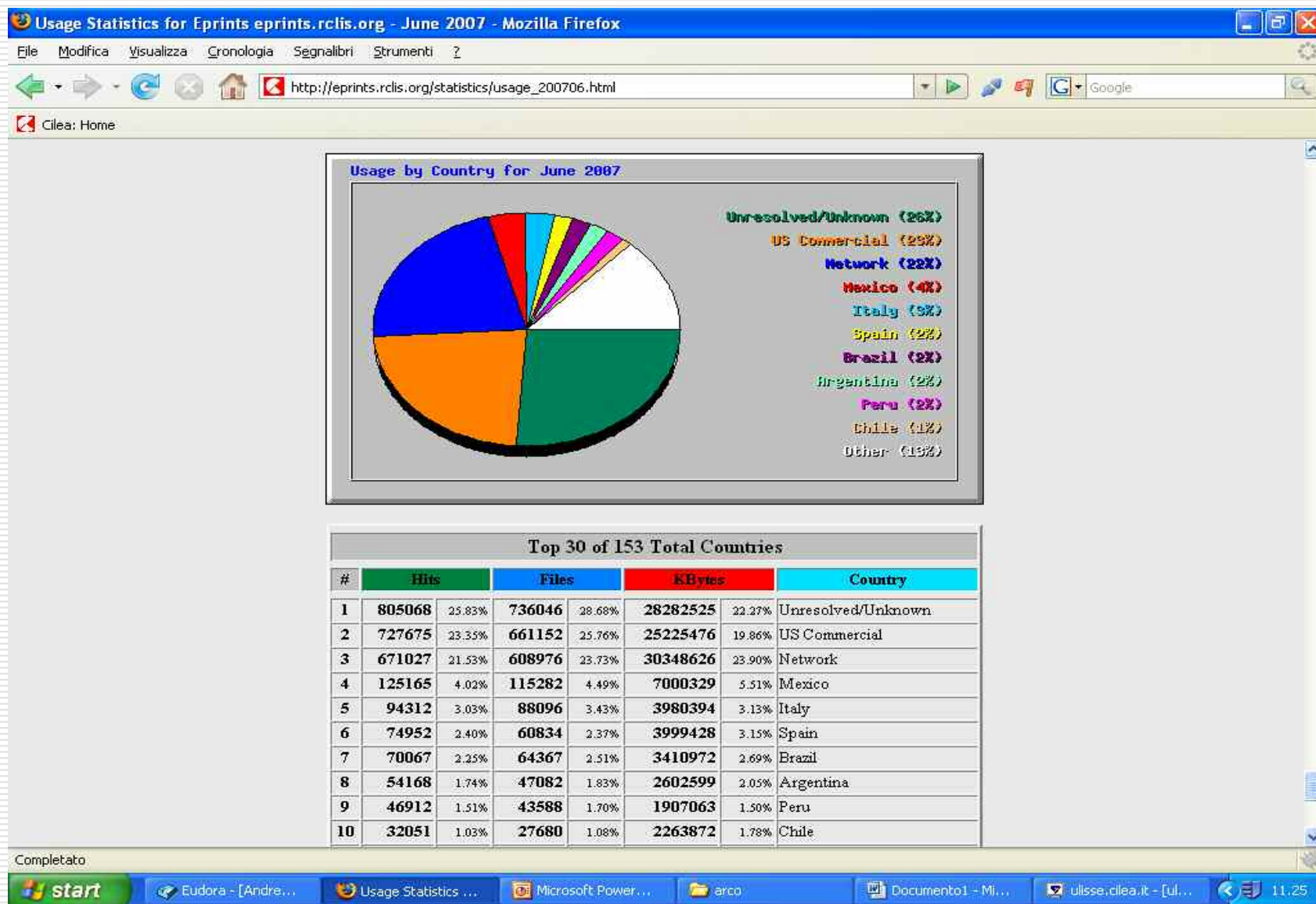
- ☐ Stessa macchina non è stessa persona
- ☐ Si vogliono vedere le visite di una persona
- ☐ Dunque quando si supera il limite massimo, uno stesso utente o un utente diverso per noi pari sono.

Come estraggo i dati.

- ☐ Non sempre dietro una macchina c'è una persona
- ☐ Ci sono software che scaricano dal web per vari motivi (robots/spiders/ agents/etc.)
- ☐ Bisogna inferire dal loro comportamento che sono software
- ☐ Una volta identificati si riconoscono successivamente usando IP e nome del software di browsing
- ☐ La lista dei robots va costantemente aggiornata
- ☐ Vi sono anche scaricatori massivi usati dalle persone, i cosiddetti "site downloader"
- ☐ Di norma sono considerati come i robots







Come estraggo i dati.

Diversi soft. su stessi dati = diverse stime

❑ **WEBALIZER:** <http://www.webalizer.org>

❑ **ANALOG:** <http://www.analog.cx/>

❑ **AWSTATS:**
<http://awstats.sourceforge.net/>

Sono software generali che vanno bene per ogni sito

Come estraggo i dati.

Localizzare tramite l'IP

- ☐ Possibile tramite studi su quanto riferito dai gestori alle autorità che gestiscono Internet
- ☐ Le opzioni gratuite sono una versione limitata di dati commerciali
- ☐ Anche questa e' una stima il cui risultato è diverso tra i vari software

Come estraggo i dati.

- ❑ AlienIP:
<http://www.iconempire.com/alien-ip/>
- ❑ GeoIP:
<http://www.maxmind.com/app/ip-location>
- ❑ Geovisite:
<http://www.geovisite.com/it/>
- ❑ HostIP.info: <http://www.hostip.info/>

Superare i limiti presentati

I cookies

- ☐ File di testo inviati dal server alla prima connessione.
- ☐ Il server può controllare sul browser la presenza o assenza del suo cookie
- ☐ Il browser può bloccarli
- ☐ L'utente può cancellarli
- ☐ Sono abbastanza temuti dagli utenti che li cancellano spesso
- ☐ Spesso usati per estrarre dati sulla navigazione in generale da agenzie pubblicitarie.
- ☐ Quest'ultima tipologia è molto avversata dagli utenti

Superare i limiti presentati

Uso delle sessioni anonime

- ☐ Usabile se invece di pagine statiche si ha un sito con pagine dinamiche
- ☐ Tengo nota dell'hand-shake con il gestionale del sito
- ☐ Scrivo in un db un identificativo di sessione e ci collego le attività rilevanti
- ☐ Chiudo la sessione quando non ho attività da quella fonte che ha fatto l'hand-shake dopo un x temporale (15-60 minuti)
- ☐ Preciso nel rilevare gli inizi
- ☐ Preciso nel rilevare l'attività
- ☐ Leggera imprecisione nella chiusura
- ☐ Non lega tra loro le diverse visite
- ☐ Non distingue persone/robots

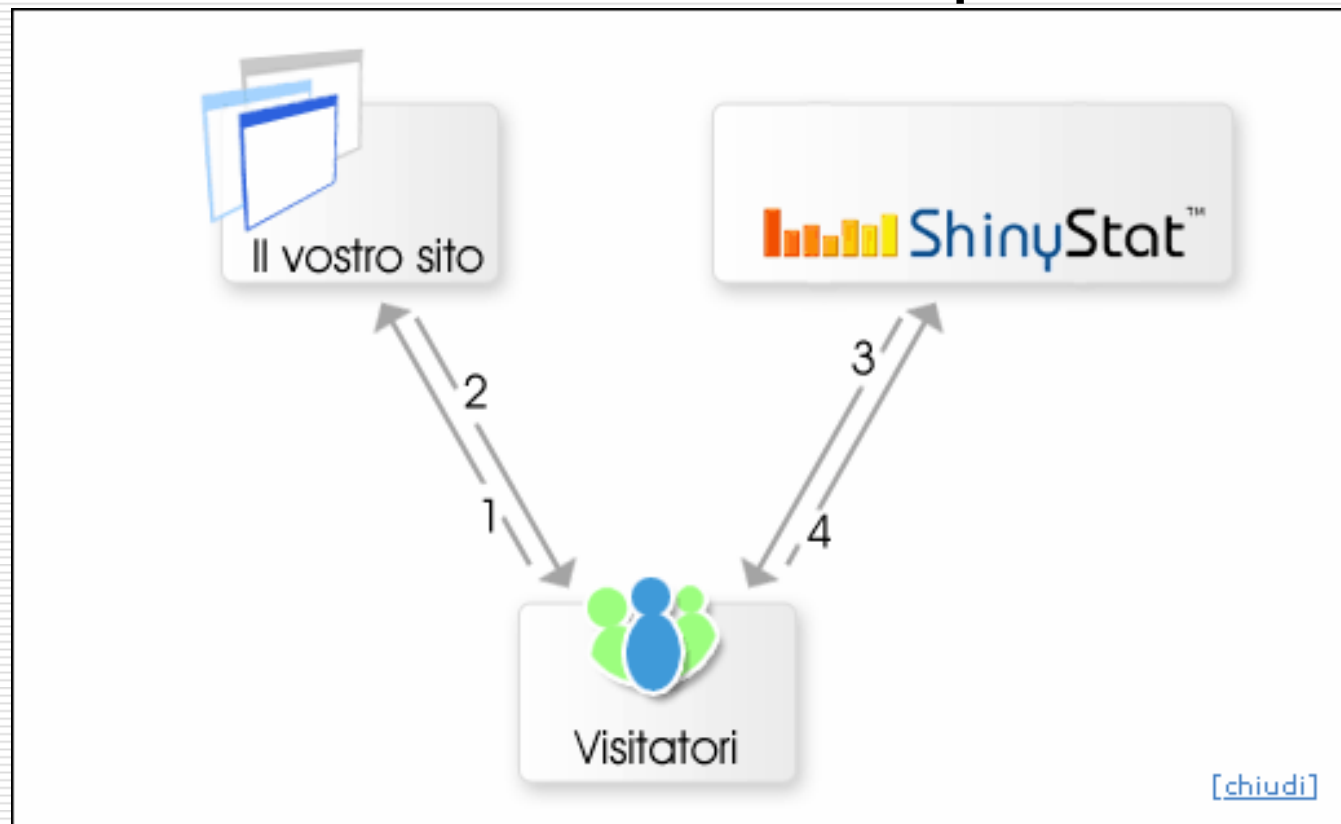
Superare i limiti presentati

Uso sessioni non anonime

- ☐ Sempre per siti dinamici
- ☐ Sempre con l'appoggio di un db
- ☐ Ogni utente e' perfettamente identificato
- ☐ Solo persone, no robots
- ☐ Le varie visite si collegano tra loro
- ☐ Le persone non amano dover ricordare una userid/password
- ☐ Voglia di anonimato diffusa
- ☐ Non si è presenti sui motori di ricerca [Elsevier ci riesce con accordi ad hoc con google scholar]
- ☐ Si esce da quello che è il web normale

Superare i limiti presentati

Portare i dati su una terza parte



Superare i limiti presentati

- ☐ Si inserisce del codice javascript in tutte le pagine
- ☐ E' il browser dell'utente che manda i dati alla terza parte
- ☐ I robots non hanno javascripts e dunque sono esclusi
- ☐ Può lasciare un traccia sul Pc dell'utente meno invadente dei cookies
- ☐ Si usa l'expertise e i mezzi di una grande organizzazione
- ☐ Posso esserci problemi di caricamento e di rete usando anche un server di terzi

Superare i limiti presentati

- ❑ Google analytics:
<http://www.google.com/analytics/>
- ❑ HiStats: <http://www.histats.com/>
- ❑ ShymyStat:
<http://www.shinystat.com>
- ❑ Site Meter:
<http://www.sitemeter.com/>

Un esempio

Le statistiche in batch di E-LIS per items

- ❑ Punto di partenza: un singolo items

- ❑ Es:

http://eprints.rclis.org/es/index.php?action=show_detail_eprint&id=6656

- ❑ I dati che si vogliono sono:

- Quanto volte sono stati letti i metadati (views)
- Quante volte sono stati scaricati i full-texts (downloads)
- Non considerare i downloads multipli
- Non considerare i robots

Un esempio

- ❑ Si opera sui logs di apache
- ❑ Troppo complesso gestire le sessioni anonime
- ❑ Si contano le righe con una specifica azione dell'applicazione
- ❑ Si fanno dei controlli sulla distanza temporale dello stesso azione fatta dallo stesso IP
- ❑ Si escludono i robots via lista fissa aggiornata periodicamente

Un esempio

- ❑ 213.58.139.xx "-" "-" "-" "GET /archive/00011300/ HTTP/1.1" 200 [...]
- ❑ Questa è l'indicazione di un view

- ❑ 213.58.139.xx "-" "-" "-" "GET /archive/00011300/01/deposito.pdf HTTP/1.1" 200 [...]
- ❑ Questa è l'indicazione di un download

- ❑ Per essere contate due volte le stesse operazioni devono essere distanziate di 180 secondi

E-LIS - Ejemplo de referenciación en bibliotecas móviles o aproximación al significado real de sus resultados - Mozilla Firefox

File Modifica Visualizza Cronologia Segnalibri Strumenti ?

http://eprints.rclis.org/archive/00011972/

Cilea: Home

E-LIS E-prints in Library and Information Science
eprints.rclis.org

The international open archive for LIS

home | about | search | browse | register | registered users area | help | FAQ | JITA

Ejemplo de referenciación en bibliotecas móviles o aproximación al significado real de sus resultados

Soto Arranz, Roberto (2007) Ejemplo de referenciación en bibliotecas móviles o aproximación al significado real de sus resultados. *Boletín de la Asociación Andaluza de Bibliotecarios*(83):pp. 11-17.

Full text available as:
PDF - Requires Adobe Acrobat Reader or other PDF viewer.

[View statistics for this eprint](#)

Abstract

[Spanish abstract]

Las bibliotecas móviles participan de una serie de factores que, si bien son imprescindibles para su funcionamiento, no pueden considerarse como estrictamente bibliotecarios. Su consideración, cuantificación y análisis son fundamentales para el correcto conocimiento de este tipo de servicios bibliotecarios. Mediante un ejemplo práctico, en este artículo se trata de poner de manifiesto esta circunstancia y reconseguir una aproximación a su verdadero significado.

[English abstract]

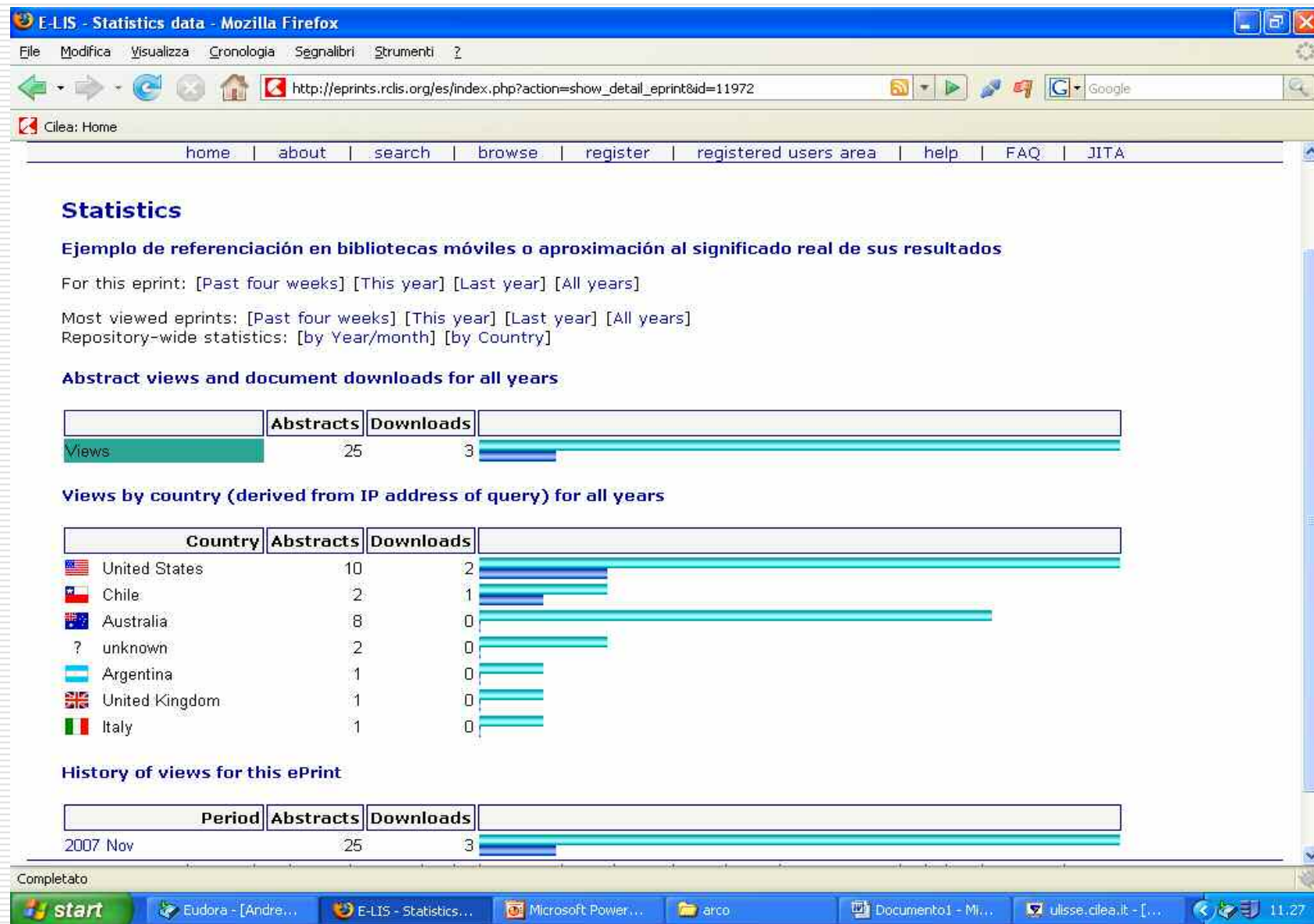
Bookmobiles have some factors that are not strictly librarians. But theses factors are very important to know their true meaning. In this text one can see a practical example about it around the comparison of three bookmobiles.

Keywords: Bibliotecas móviles, bibliobuses, referenciación, evaluación, Mobiles libraries, bookmobiles, comparison, evaluation

Subjects: F. Management. > FZ. None of these, but in this section.

Completo

start Eudora - [Andre... E-LIS - Ejemplo ... Microsoft Power... arco Documento1 - Mi... ulisse.cilea.it - [... 11.26



Un esempio

Punti problematici

- ❑ La geolocalizzazione
- ❑ La lista dei robots non si aggiorna automaticamente
- ❑ Notevole spazio necessario sul server
- ❑ Soluzione non confrontabile con indicatori di utilizzo sito standard, dunque compresenza di statistiche dovute a webanalyzer (a partire da <http://eprints.rclis.org/statistics-sheet.html>)

DOMANDE ?

Link utili

❑ "Informatica di base", R. Gaeta, 2004
[URL: http://www.di.unito.it/~rossano/DIDATTICA/INF-0304/#Lucidi](http://www.di.unito.it/~rossano/DIDATTICA/INF-0304/#Lucidi)

❑ modulo Apache 2 mod_log_config:
http://httpd.apache.org/docs/2.0/mod/mod_log_config.html

RFC 2616:
<http://www.w3.org/Protocols/rfc2616/rfc2616.html>

❑ Sul dynamic DNS:
http://it.wikipedia.org/wiki/Dynamic_DNS